

## PATENT ABSTRACTS OF JAPAN

(11)Publication number : 2000-231395

(43)Date of publication of application : 22.08.2000

(51)Int.Cl. G10L 13/06  
G10L 13/08

(21)Application number : 11-030684

(71)Applicant : NIPPON TELEGR &amp; TELEPH CORP &lt;NTT&gt;

(22)Date of filing : 08.02.1999

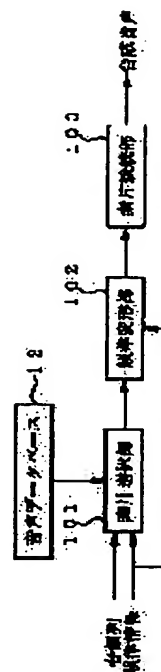
(72)Inventor : MIZUNO HIDEYUKI  
TANAKA KIMITO  
NAKAJIMA SHINYA  
ABE MASANOBU

## (54) METHOD AND DEVICE FOR SYNTHESIZING VOICE

## (57)Abstract:

**PROBLEM TO BE SOLVED:** To provide method and device for synthesizing a voice scalably changing the voice data and synthetic quality according to a use from a synthesized voice of quality of an extent obtained with a voice synthetic method based on the voice data of low capacity using a voice elemental piece of a length of a phoneme or a syllable etc., until the synthesized voice of high quality similar to a natural voice based on a voice data base of large capacity.

**SOLUTION:** When a phoneme line and rhythm information are inputted to an elemental piece selection part 101, the elemental piece selection part 101 refers to the inputted phoneme line and rhythm information, and selects the optimum voice elemental piece data from the voice data base 12 to send them to a rhythm deformation part 102. A voice waveform, the rhythm information and rhythm boundary information, etc., are stored in the voice data base 12. The rhythm deformation part 102 deforms the voice elemental piece data selected by the elemental piece selection part 101 so as to be suited to the inputted rhythm information to send them to an elemental piece connection part 103. The elemental piece connection part 103 connects successively the elemental piece data deformed by the rhythm deformation part 102 to generate the synthesized voice.



(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開2000-231395

(P2000-231395A)

(43) 公開日 平成12年8月22日 (2000.8.22)

(51) Int.Cl.<sup>7</sup>

識別記号

F I

テマコード (参考)

G 1 0 L 13/06

G 1 0 L 5/04

F 5 D 0 4 5

13/08

3/00

H 9 A 0 0 1

審査請求 未請求 請求項の数6 O L (全 8 頁)

(21) 出願番号

特願平11-30684

(22) 出願日

平成11年2月8日 (1999.2.8)

(71) 出願人 000004226

日本電信電話株式会社

東京都千代田区大手町二丁目3番1号

(72) 発明者 水野 秀之

東京都新宿区西新宿三丁目19番2号 日本

電信電話株式会社内

(72) 発明者 田中 公人

東京都新宿区西新宿三丁目19番2号 日本

電信電話株式会社内

(74) 代理人 100064908

弁理士 志賀 正武

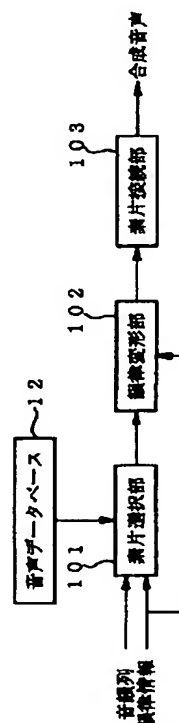
最終頁に続く

(54) 【発明の名称】 音声合成方法及び装置

(57) 【要約】

【課題】 音素または音節等の長さの音声素片を使用した低容量の音声データに基づく音声合成方法で得られる程度の品質の合成音声から、大容量の音声データベースに基づく自然音声と同様の高品質な合成音声まで、用途に応じてスケーラブルに音声データと合成品質を変更可能な音声合成方法及び装置を提供する。

【解決手段】 音韻列と韻律情報が素片選択部101に入力されると、素片選択部101は入力された音韻列と韻律情報を参照して音声データベース12より最適な音声素片データを選択して韻律変形部102に送る。この音声データベース12には音声波形、音韻情報、音韻境界情報などが格納されている。韻律変形部102は、入力された韻律情報に適合するように素片選択部101で選択された音声素片データを変形して素片接続部103に送る。素片接続部103は韻律変形部102で変形された素片データを順に接続して合成音声を生成する。



## 【特許請求の範囲】

【請求項1】 入力された音韻列と韻律情報に対応づけられた音声素片データを音声データベースから選択して順次接続することにより音声信号を合成する音声合成方法において、

前記入力された音韻列を予め決められた規則に従って部分音韻列に分解する分解過程と、

対応づけられた音韻列が前記分解された部分音韻列と一致し、かつ該音韻列の前後の音韻が前記部分音韻列の前後の音韻と一致する音声素片データの前記音声データベース中における存在の有無を判断する判断過程と、

前記音声素片データが存在する場合には、当該音声素片データを選択する選択過程と、

前記音声素片データが存在しない場合には、前記分解された部分音韻列を前記入力された音韻列として、前記部分音韻列の長さが予め定められた最小音韻長に分解されるまで前記分解過程と前記判断過程を反復させる過程と、

前記部分音韻列に対応する前記韻律情報を構成する部分韻律情報に応じて、前記選択された音声素片データを韻律変形する過程と、

前記韻律変形を受けた音声素片データを順次接続して音声信号を合成する過程とを有することを特徴とする音声合成方法。

【請求項2】 前記最小音韻長にまで分解された部分音韻列が存在せず、かつ前記最小音韻長が2である場合、前記部分音韻列と該部分音韻列の前後の音韻を含む部分音韻列を連鎖音韻に分解する過程をさらに有し、前記選択過程では前記連鎖音韻に対応する音声素片データを選択することを特徴とする請求項1記載の音声合成方法。

【請求項3】 一個の部分音韻列について前記選択された音声素片データが複数個存在する場合、それら音声素片データに対応する韻律と前記部分音韻列に対応する韻律との類似性を判断する過程をさらに有し、前記選択過程では、前記複数個の音声素片データのうち、最も類似性の高い音声素片データを選択することを特徴とする請求項1乃至2記載の音声合成方法。

【請求項4】 入力された音韻列と韻律情報に対応づけられた音声素片データを音声データベースから選択して順次接続することにより音声信号を合成する音声合成装置において、

前記入力された音韻列を予め決められた規則に従って部分音韻列に分解する分解手段と、

対応づけられた音韻列が前記分解された部分音韻列と一致し、かつ該音韻列の前後の音韻が前記部分音韻列の前後の音韻と一致する音声素片データの前記音声データベース中における存在の有無を判断する判断手段と、

前記音声素片データが存在することを条件として、当該音声素片データを選択する選択手段と、

前記音声素片データが存在しないことを条件として、前記分解された部分音韻列を前記入力された音韻列として前記分解手段に入力して、前記部分音韻列の長さが予め定められた最小音韻長に分解されるまで前記分解手段と前記判断手段とを反復動作させるように制御する手段と、

前記部分音韻列に対応する前記韻律情報を構成する部分韻律情報に応じて、前記選択された音声素片データを韻律変形する手段と、

10 前記韻律変形を受けた音声素片データを順次接続して前記音声信号を合成する手段とを具備することを特徴とする音声合成装置。

【請求項5】 前記最小音韻長にまで分解された部分音韻列が存在せず、かつ前記最小音韻長が2であることを条件として、前記部分音韻列と該部分音韻列の前後の音韻を含む部分音韻列を連鎖音韻に分解する手段をさらに具備し、

前記選択手段は前記連鎖音韻に対応する音声素片データを選択することを特徴とする請求項4記載の音声合成装置。

【請求項6】 一個の部分音韻列について前記選択された音声素片データが複数個存在することを条件として、それら音声素片データに対応する韻律と前記部分音韻列に対応する韻律との類似性を判断する手段をさらに具備し、

前記選択手段は、前記複数個の音声素片データのうち、最も類似性の高い音声素片データを選択することを特徴とする請求項4乃至5記載の音声合成装置。

## 【発明の詳細な説明】

## 【0001】

【発明の属する技術分野】この発明は、テキストを入力しそのテキストに応じた任意の音声を合成する音声合成方法及び装置に関し、特に、主に音韻列と韻律情報とから音声合成する規則音声合成方法及びこの方法を実現するための装置に関するものである。

## 【0002】

【従来の技術】従来の音声合成方法では、あらかじめ、音声素片として音素単位や、CV、VCV、CVC

(C：子音、V：母音)など音韻の調音結合を考慮した単位、3音韻以上のフォルマントを考慮した単位、または前記全ての単位で音声データベースを作成しておき、音声合成の際に、入力テキストや韻律情報に応じて音声データベース中から適切な素片データを選択して接続することによって音声合成を行っているものが多い

(特開昭59-204097号公報、特開平1-078300号公報、特開平6-095692号公報、特開平9-090972号公報)。この音声合成方法では、合成音声の品質はおおよそ使用する音声データベースの容量と比例しており、容量は少ないが自然音声よりかなり劣ったものから、容量は大きいがある程度高品質なもの

まで様々なものが開発・製品化されている。しかし、それらの製品は全く独立に開発されており互換性等がないため、容量、品質、応答時間などの使用条件に応じて使い分けることが困難である。

【0003】さらに、近年では大容量な記憶装置の使用コストの低下にともなって、数十分から数時間に及ぶ音声データをそのまま大容量の記憶装置に蓄積し、入力されたテキスト及び韻律情報に応じた適当な基準で大容量の音声データから適当な長さの音声素片を切り出すとともに、入力された韻律情報に従って切り出された音声素片を適切に変形し接続することによって合成する音声合成方法も提案されている(特許第2761552号)。この方法では大容量の音声データを用意することで、理論的には高品質な合成音声を作成することが可能であるが、大容量の音声データとそれを格納する記憶装置が必要であるためシステム価格が高くなることや、音声データを収集する基準または方法が確立されていないため、必要な品質に見合った最適な規模の音声データを収集することが不可能であること、存在する音声データから適切な音声素片を切り出す最適な規則や方法が確立されていないため、切り出された音声素片が必ずしも適切でなく合成音声全体の品質が安定しないこと等の問題がある。

#### 【0004】

【発明が解決しようとする課題】この発明は上述した問題点に鑑みてなされたものであり、その目的は、音素または音節等の長さの音声素片を使用した低容量の音声データに基づく音声合成方法によって得られる程度の品質の合成音声から、大容量の音声データベースに基づく自然音声と同様の高品質な合成音声まで、用途に応じてスケラブルに音声データと合成品質を変更することが可能な音声合成方法及び装置を提供することにある。また、この発明の目的は、大容量の音声データにもとづく音声合成方式の問題を解決し、音声データの収集基準と音声素片の選択規則を明確化することにより、常に適切な音声素片データの選択が保証された高品質な合成音声を実現できる音声合成方法及び装置を提供することにある。

#### 【0005】

【課題を解決するための手段】以上の課題を解決するために、請求項1記載の発明は、入力された音韻列と韻律情報に対応づけられた音声素片データを音声データベースから選択して順次接続することにより音声信号を合成する音声合成方法において、前記入力された音韻列を予め決められた規則に従って部分音韻列に分解する分解過程と、対応づけられた音韻列が前記分解された部分音韻列と一致し、かつ該音韻列の前後の音韻が前記部分音韻列の前後の音韻と一致する音声素片データの前記音声データベース中における存在の有無を判断する判断過程と、前記音声素片データが存在する場合には、当該音声

素片データを選択する選択過程と、前記音声素片データが存在しない場合には、前記分解された部分音韻列を前記入力された音韻列として、前記部分音韻列の長さが予め定められた最小音韻長に分解されるまで前記分解過程と前記判断過程を反復させる過程と、前記部分音韻列に対応する前記韻律情報を構成する部分韻律情報に応じて、前記選択された音声素片データを韻律変形する過程と、前記韻律変形を受けた音声素片データを順次接続して音声信号を合成する過程とを有することを特徴としている。

【0006】また、請求項2記載の発明は、請求項1記載の発明において、前記最小音韻長にまで分解された部分音韻列が存在せず、かつ前記最小音韻長が2である場合、前記部分音韻列と該部分音韻列の前後の音韻を含む部分音韻列を連鎖音韻に分解する過程をさらに有し、前記選択過程では前記連鎖音韻に対応する音声素片データを選択することを特徴としている。また、請求項3記載の発明は、請求項1乃至2記載の発明において、一個の部分音韻列について前記選択された音声素片データが複数個存在する場合、それら音声素片データに対応する韻律と前記部分音韻列に対応する韻律との類似性を判断する過程をさらに有し、前記選択過程では、前記複数個の音声素片データのうち、最も類似性の高い音声素片データを選択することを特徴としている。

【0007】また、請求項4記載の発明は、入力された音韻列と韻律情報に対応づけられた音声素片データを音声データベースから選択して順次接続することにより音声信号を合成する音声合成装置において、前記入力された音韻列を予め決められた規則に従って部分音韻列に分解する分解手段と、対応づけられた音韻列が前記分解された部分音韻列と一致し、かつ該音韻列の前後の音韻が前記部分音韻列の前後の音韻と一致する音声素片データの前記音声データベース中における存在の有無を判断する判断手段と、前記音声素片データが存在することを条件として、当該音声素片データを選択する選択手段と、前記音声素片データが存在しないことを条件として、前記分解された部分音韻列を前記入力された音韻列として前記分解手段に入力して、前記部分音韻列の長さが予め定められた最小音韻長に分解されるまで前記分解手段と前記判断手段とを反復動作させるように制御する手段と、前記部分音韻列に対応する前記韻律情報を構成する部分韻律情報に応じて、前記選択された音声素片データを韻律変形する手段と、前記韻律変形を受けた音声素片データを順次接続して前記音声信号を合成する手段とを具備することを特徴としている。

【0008】また、請求項5記載の発明は、請求項4記載の発明において、前記最小音韻長にまで分解された部分音韻列が存在せず、かつ前記最小音韻長が2であることを条件として、前記部分音韻列と該部分音韻列の前後の音韻を含む部分音韻列を連鎖音韻に分解する手段をさ

らに具備し、前記選択手段は前記連鎖音韻に対応する音声素片データを選択することを特徴としている。また、請求項6記載の発明は、請求項4乃至5記載の発明において、一個の部分音韻列について前記選択された音声素片データが複数個存在することを条件として、それら音声素片データに対応する韻律と前記部分音韻列に対応する韻律との類似性を判断する手段をさらに具備し、前記選択手段は、前記複数個の音声素片データのうち、最も類似性の高い音声素片データを選択することを特徴としている。

【0009】以上のように、本発明は、音韻情報と韻律情報とから音声合成する規則音声合成方法及び装置に適用されるものである。そして本発明は、入力された音韻情報に従って音声データベースから音声素片データを選択する際に、音素や音節などの一定の単位での選択または複雑な規則や計算に基づく選択を行うのではなく、音韻情報にある単純な規則に従って部分音韻列に分解し、分解された音韻列およびその前後の音韻環境に適合する音声素片データを音声データベースから選択し、適合する音声素片データが無かった部分音韻列のみをさらに別の単純な規則に従って分解し、その分解された音韻列に適合する音声素片データを音声データベースから選択し、さらに適合する音声素片データが無かった部分音韻列のみ分解し、という多段階の分解と選択を入力音韻列に対応する全ての音声素片データが見つかるまで行うことに特徴を有している。

【0010】このように多段階の選択を行うことで、最下段の分解規則に対応した最小単位で音声データベースを構成した場合が最も低容量・低品質な用途に対応するとともに、それより上の段階の分解規則に対応した単位の音声素片データを音声データベースに追加することで、より高品質な用途に対応させることが可能となり、また、最上段の分解規則に対応した最長単位の音声素片データが全て音声データベースに存在する場合は最高品質の用途に対応させることが可能となる。そして各段階に対応する音声素片データを音声データベースに追加または削除するだけで、音声合成システムの変更が簡単に実現できる。

【0011】また、本発明では最終段階の音声素片選択においては、環境を考慮しない連鎖音韻(CV, VV, VC)にもとづく合成方法も適用可能であることに特徴を有している。このようにすることで、環境を考慮して音韻単位で音声素片データを用意した場合は数千~数万個の音声素片データが必要となるのに対し、本発明のように連鎖音韻単位で音声素片データを用意した場合は約千個程度の音声素片データを用意すればよい。そのため少量の記憶装置やメモリなどに音声データベースを格納でき、LSI(大規模集積回路)への内蔵用途等にも対応可能となる。

【0012】また、本発明では各段階で複数の音声素片

データが選択された場合、音声素片のピッチパターンが合成すべきピッチパターンともっとも類似する音声素片データを選択することにも特徴を有している。このようにすることで、音声合成時のピッチの変更量を少なくすることができ、合成音声の品質を向上させることが可能となる。また合成すべきピッチパターンと同一のピッチデータをもつ音声素片データを追加することで、ピッチの変形処理が不要になり、その場合の品質は編集音声合成の品質とほぼ同等となる。

# 10 【0013】

【発明の実施の形態】以下、図面を参照してこの発明の一実施形態を述べる。図1に本実施形態による音声合成処理を実現するための音声合成装置の基本構成を示す。図示したように、この音声合成装置は音声データベース12、素片選択処理を行う素片選択部101、韻律変形処理を行う韻律変形部102、素片接続処理を行う素片接続部103から構成されている。これら各部が行うそれぞれの処理については以下に詳述する。また、図2はこの音声合成処理の手順を示したフローチャートである。最初に図1を参照しながら音声合成装置の全体動作について説明し、その後、図3及び図4を参照して音声合成装置を構成する幾つかの機能ブロックの構成及びその動作の詳細について説明する。

【0014】図1に示すように、音韻列と韻律情報が素片選択部101に入力される(図2のステップS1)と、素片選択部101は入力された音韻列と韻律情報を参照して音声データベース12より最適な波形(素片データ)を選択して韻律変形部102に送る(ステップS2)。ここで、音声データベース12には音声波形、音韻情報、音韻境界情報などが格納されているものとする。なお、素片選択部101の詳細な構成については後述する。次に、韻律変形部102は、入力された韻律情報に適合するように、部分音韻列に対応する韻律情報を構成している部分韻律情報に応じて、前記素片選択部101で選択された素片データを変形して素片接続部103に送る(ステップS3)。

【0015】ここで、波形の変形方法としては、PSOLA法(E.Moulines and F.Charpentier, "Pitch-synchronous waveform processing techniques for text-to-speech synthesis using diphones", Speech Communication, Vol.9, pp.453-467, 1990.12)、IPSE法(田中ら, 「基本周波数に応じてスペクトル包絡を変形するテキスト合成システム」, 信学技報, SP96-130, pp.23-30, 1997.3)、STRAIGHT法(河原, 「聴覚の情景分析と高品質音声分析変換合成法STRAIGHT」, 音響学会講演論文集, pp.189-192, 1997.9)などがある。従って、音声データベース12には、それらの合成方式に応じて最適な形で格納すればよく、必ずしも波形データをそのまま格納する必要はない。例えばSTRAIGHT法を用いるのであれば、事前にSTRA

I G H T分析で得られたパラメータを格納しておくことで音声合成時の計算時間が削減できる。そして最後に、素片接続部103は、前記韻律変形部102で変形された素片データを順に接続し合成音声を生成する(ステップS4)。以上が本実施形態による音声合成装置において行われる処理の全体的な流れである。

【0016】次に、素片選択部101における処理の1例について、図3のブロック図と前掲した図2のフローチャートを参照して説明する。なお、図3において図1に示したものと同一構成要素については同一の符号を付してある。図3に示したように、素片選択部101は第一段階分解部201、第一段階選択部202、韻律マッチング部203、第二段階分解部204、第二段階選択部205、第三段階分解部206、第三段階選択部207、最終段階選択部208から構成されている。まず、例えば入力音韻列を“bakuoNga/giNsekaino”(ここで、記号/はアクセント境界を示す)とした場合、第一段階分解部201で入力音韻列を例えばアクセント句単位に分する(ステップS21)。これは、日本語ではアクセント単位でまとまって発声される場合が多く、アクセント句が発声現象の大きなまとまりと考えられるためである。この例では、前記入力音韻列が部分音韻列“bakuoNga”とgiNsekaino”に分解される。

【0017】次に、第一段階選択部202は音声データベース12から音韻列“bakuoNga”, 前音韻環境が語頭(図中の記号#), 後音韻環境が“g”という素片データ、および、音韻列“giNsekaino”, 前音韻環境が“a”, 後音韻環境が語尾(図中の記号#)という素片データを音声データベース12からそれぞれ検索する(ステップS22)。図で示すとおり、音韻列“giNsekaino”に対応する素片データが見つからず(同ステップが“NO”)に、音韻列“bakuoNga”に対応する素片データ21のみ見つかった(同ステップが“YES”)場合、第一段階選択部202は素片データ21のみを韻律マッチング部203に送る。この場合、音韻列“bakuoNga”に対応する素片データは1つしかない(ステップS25が“NO”)ため、韻律マッチング部203は素片データ21をそのまま図1の韻律変形部102に送る。また、第一段階選択部202は音韻列“giNsekaino”を第二段階分解部204に送る(ステップS24)。

【0018】次に、第二段階分解部204は、前記第一段階選択部202による音声データベース12の検索で見つからなかった音韻列“giNsekaino”を例えば音節に母音や撥音の連続を含む単位で分解する(ステップS21)。これは撥音や母音が連続している場合、発声現象的に連続しており音響的にも境界を設定するのが困難であるためである。そしてこの例では、“giN”; “se”, “kai”, “no”の4つの部分音韻列に分解される。次に、第二段階選択部205は第一段階選択部202と同様に、音声データベース12から、音韻列“gi

N”, 前音韻環境が“a”, 後音韻環境が“s”の素片データと、音韻列“se”, 前音韻環境が“N”, 後音韻環境が“k”の素片データと、音韻列“kai”, 前音韻環境が“e”, 後音韻環境が“n”の素片データと、音韻列“no”, 前音韻環境が“i”, 後音韻環境が語尾(図中の記号#)の素片データをそれぞれ検索する(ステップS22)。

【0019】この結果、図で示すとおり第二段階選択部205は“giN”に対応する素片データ22及び素片データ23, “se”に対応する素片データ24, “no”に対応する素片データ25を韻律マッチング部203に送る(ステップS23が“YES”)。この場合、素片データ24と素片データ25はいずれも音韻列に対応する素片が1つである(ステップS25が“NO”)ため、韻律マッチング部203は素片データ21と同様にこれらをそのまま図1の韻律変形部102に送る。一方、素片データ22と素片データ23(ステップS25が“YES”)については、韻律マッチング部203が入力された韻律情報とマッチングを行い、入力韻律情報と最も近い素片データを選択してから図1の韻律変形部102に送る(ステップS26)。

【0020】ここで、韻律の近さの判定方法は使用する音声データベース12の構成による。例えば、音声データベース12がピッチのバリエーションについてのみ考慮した音声データベースであれば、入力ピッチパターンと最も近い(最も類似性の高い)ピッチパターンをもつ素片データを選ぶことで十分である。また、特に韻律等を考慮していない音声データベースを使用するのであれば、平均ピッチ、ピッチ形状、時間長、パワーの各韻律パラメータについて、入力された値と素片データの持つ値との差分の絶対値を求め、これら絶対値に対して各韻律パラメータ毎の重み係数を掛けて足し合わせることで韻律コストを求め、その値の小さいものを選ぶことが望ましいと考えられる(広川ら, “波形編集型規則合成法における波形選択関数の検討”, 音響学会講演論文集, pp.157-158, 1989.3)。この例では、素片データ22が入力ピッチパターンに近いと判断されたとして、韻律マッチング部203は素片データ22を図1の韻律変形部102に送る。

【0021】次に、第三段階分解部206は、前記第二段階選択部205による音声データベース12の検索で見つからなかった(ステップS23が“NO”, ステップS24)部分音韻列“kai”を例えば音節に分解する(ステップS21)。これは、音節の構成要素である子音と母音は音響的にも発声現象的にも密接に結びついているため、分離して取り扱うのは音質の劣化を招く可能性が大きいためである。この例では、部分音韻列“kai”が“ka”と“i”の2つの部分音韻列に分解される。次に、第三段階選択部207は第一段階選択部202及び第二段階選択部205と同様に、音声データベ



ス12から、音韻列"ka"、前音韻環境が"e"、後音韻環境が"i"の素片データと、音韻列"i"、前音韻環境が"a"、後音韻環境が"n"の素片データをそれぞれ検索する(ステップS22)。

【0022】図で示すとおり音韻列"i"に対応する素片データ26が一つ見つかり(ステップS23が"YES")、音韻列"ka"に対応する素片データが見つからなかった(ステップS23が"NO")とする。すると第三段階選択部207は、韻律マッチング部203に素片データ26を送り、韻律マッチング部203は音韻列"i"に対応する素片が1つだけ(ステップS25が"NO")のため、前記同様に素片データ26を図1の韻律変形部102に送る。最後に、最終段選択部208は前記第三段階選択部207による音声データベース12の検索で見つからなかった部分音韻列"ka"を選択する(ステップS24、ステップS21~S26)。

【0023】次に、最終段選択部208の詳細について図4を参照して以下に説明する。なお、図4において図1又は図3に示したものと同一構成要素については同一の符号を付してある。この図4には2種類の分解・選択方法の一例について示してある。図4(a)では、前記の第一から第三までの分解処理および選択処理と同様に、前記第三段階選択部207による音声データベース12の検索で見つからなかった部分音韻列"ka"を音声の基本単位である音韻に分解する方法の一例について示してある。この例では音韻分解部301にて部分音韻列"ka"が"k"と"a"に分解される。

【0024】次に、音韻選択部302では第一段階選択部202、第二段階選択部205と同様に、音声データベース12から、音韻"k"、前音韻環境が"e"、後音韻環境が"a"の素片データと、音韻列"a"、前音韻環境が"k"、後音韻環境が"i"の素片データをそれぞれ検索する。図4(a)で示すとおり"k"に対応する素片データ31および"a"に対応する素片データ32が見つかったとして、音韻選択部302は図3の韻律マッチング部203に素片データ31と素片データ32を送る。韻律マッチング部203はこれら各音韻に対応する素片が各々1つのみのため、前記同様にこれら素片データを図1の韻律変形部102に送る。なお、もしも対応する素片データが複数あった場合は、前記同様に韻律マッチング部203にて入力された韻律情報と最も近い韻律を持つ素片データを選択しそれを韻律変形部102に送る。

【0025】一方、図4(b)では、前記第三段階選択部207による音声データベース12の検索で見つからなかった部分音韻列"ka"を連鎖音韻に分解する方法の一例について示してある。連鎖音韻に分解する理由は、前記図4(a)で示す方法と比較した場合に、少ないデータ量であらゆる音韻列の合成が可能となるためである。前記図4(a)に基づく方法では数千~数万の素片

データが必要なのに対し、図4(b)で示す方法では約1000個程度の素片データのみでよいため、より少ない記憶容量で音声合成が実現可能となる。この例では、部分音韻列"ka"の前音韻環境が"e"、後音韻環境が"i"であったことから、連鎖音韻分解部303はこれを"ek"、"ka"、"ai"に分解して連鎖音韻選択部304に送る。

【0026】次に、連鎖音韻選択部304は音声データベース12から"ek"、"ka"、"ai"である連鎖音韻の素片データをそれぞれ検索する。図4(b)で示すとおり"ek"に対応する素片データ33、"ka"に対応する素片データ34、"ai"に対応する素片データ35が見つかったとして、連鎖音韻選択部304は図3の韻律マッチング部203に素片データ33~素片データ35を送る。この場合、各連鎖音韻に対応する素片データは各々1つのみのため、韻律マッチング部203は前記同様にこれら素片データを図1の韻律変形部102に送る。なお、もしも対応する素片データが複数ある場合は、前記同様に韻律マッチング部203にて入力韻律情報と最も近い韻律を持つ素片データを選択しそれを韻律変形部102に送る。またこの後、素片接続部103が韻律変形部102で変形された素片データ(ステップS3)を順に接続して合成音声を作成する(ステップS4)が、図4(b)の場合は"k"と"a"が重複するので、このまま素片データを接続するだけでは音韻の重複が避けられない。そのため、素片接続部103は素片接続に先立って重複しないように音韻の中間部分をつなぐようにしている。

【0027】

【発明の効果】以上述べたように、この発明によれば、入力音韻列と入力韻律情報に対して段階的に音声データベースから音声素片データを選択してそれら音声素片データを接続することで出力音声を合成している。このため、低コストの合成システムから高コストではあるが高品質な合成システムまで用途に応じてシステム規模をスケラブルに変更可能であり、実用性に優れた合成システムの提供が可能である。また、最低段階の選択規則に対応した音声素片データに基づく合成品質は保証されているため、一定以上の品質が保証された合成音声の提供が可能である。

【0028】また、請求項2又は5記載の発明では、部分音韻列を前後の音韻環境を含めて連鎖音韻に分解して連鎖音韻単位で音声素片データを選択するようにしている。このため、環境を考慮して音韻単位で音声素片データを用意した場合には数千~数万個の音声素片データが必要となるのに対し、連鎖音韻単位で音声素片データを用意することで約千個程度の音声素片データを用意すれば良くなる。そのため、音声データベースを少容量のメモリ等に格納することができ、LSIに内蔵するなどの用途にも対応することができる。

【0029】また、請求項3又は6記載の発明では、ある部分音韻列について複数の音声素片データが選択された場合に、例えば、音声素片データのピッチパターンが合成すべきピッチパターンに最も類似するものを選択している。このため、例えば音声合成時のピッチの変更量が少なくなつて合成音声の品質を向上できるほか、合成すべきピッチパターンと同一のピッチデータを持つ音声素片データを追加すればピッチの変形処理が不要になるため、編集音声合成の品質とほぼ同等の合成音声を得られる。

【図面の簡単な説明】

【図1】 本発明の一実施形態による音声合成装置の構成を示すブロック図である。

【図2】 同実施形態における音声合成方法の手順を示したフローチャートである。

【図3】 図1に示す素片選択部101の詳細な構成を示すブロック図である。

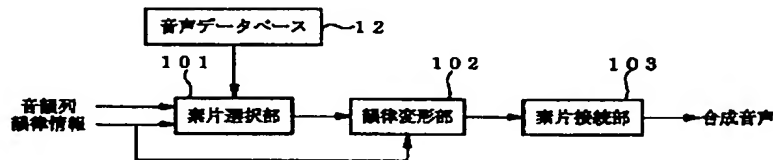
【図4】 図3に示す最終段階選択部208の詳細な構成を示すブロック図であつて、(a)は部分音韻列を音韻に分解するようにした場合の構成例、(b)は部分音韻列を連鎖音韻に分解するようにした場合の構成例であ

＊る。

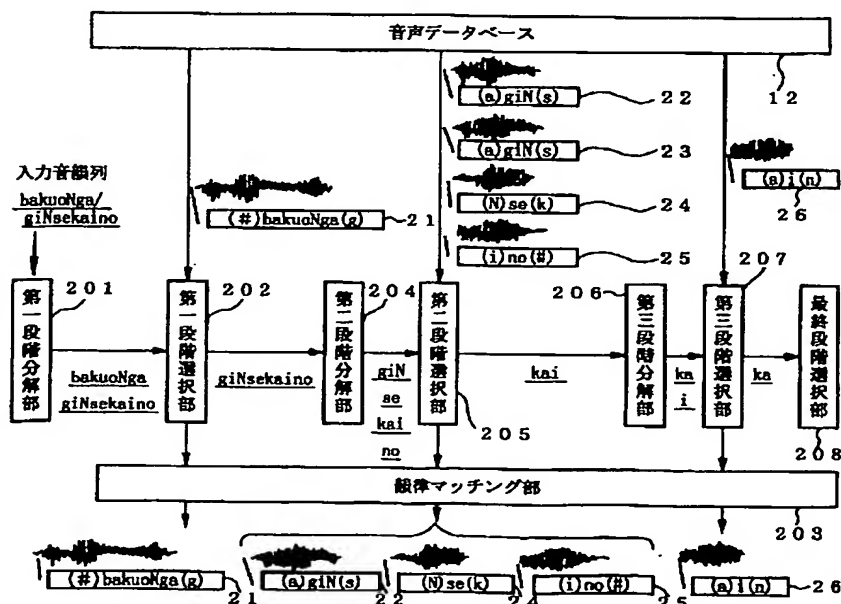
【符号の説明】

- 12 音声データベース
- 21～26, 31～35 素片データ
- 101 素片選択部
- 102 韻律変形部
- 103 素片接続部
- 201 第一段階分解部
- 202 第一段階選択部
- 203 韻律マッチング部
- 204 第二段階分解部
- 205 第二段階選択部
- 206 第三段階分解部
- 207 第三段階選択部
- 208 最終段階選択部
- 301 音韻分解部
- 302 音韻選択部
- 303 連鎖音韻分解部
- 304 連鎖音韻選択部

【図1】

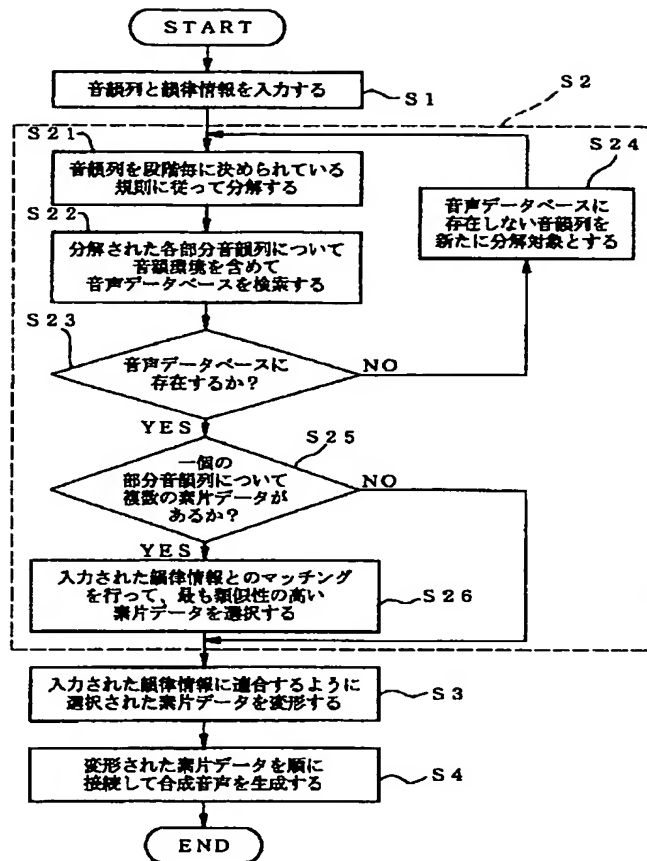


【図3】

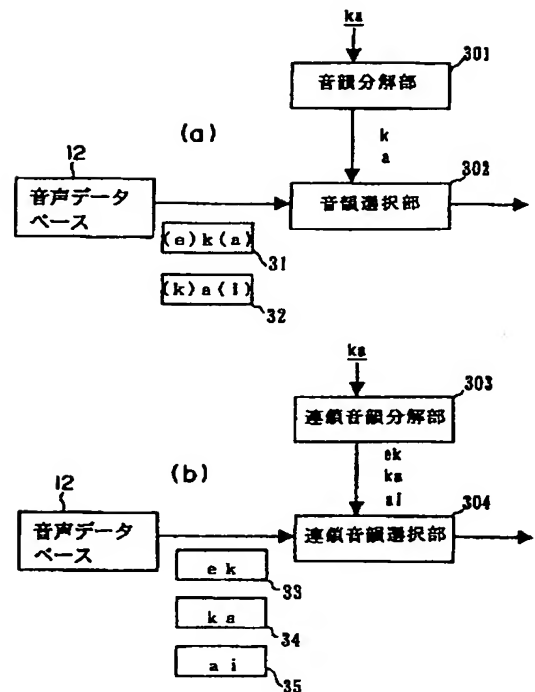




【図2】



【図4】



フロントページの続き

(72)発明者 中島 信弥  
東京都新宿区西新宿三丁目19番2号 日本  
電信電話株式会社内

(72)発明者 阿部 匡伸  
東京都新宿区西新宿三丁目19番2号 日本  
電信電話株式会社内  
Fターム(参考) 5D045 AA07 AB01  
9A001 HH18 JJ01